

Generalized Linear Models

Lecture 10: Nonparametric regression



$$y_i = f(x_i) + \varepsilon_i$$

- How to estimate f?
- Could assume $f(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$
- Or $f(x) = \beta_0 + \beta_1 x^{\beta_2}$
- OK if you know the right form.
- But often better to assume only that *f* is continuous and smooth.

Examples



1 Kernel estimators

2 Local polynomials





In brief

For any point x_0 , the value of the function at that point $f(x_0)$ is some combination of the (nearby) observations.

Kernel function

The contribution of each observation x_j , $f(x_j)$ to $f(x_0)$ is calculated using a weighting function or Kernel $K_h(x_0, x_j)$. h is the width of the neighborhood.

A simple estimate of $f(x_0)$ at any point x_0 is the mean of the k points closest to x_0 .

$$\hat{f}(x_0) = \operatorname{Ave}(y_i | x_i \in N_k(x_0)).$$



True function KNN average Observations contributing to $\hat{f}(x_0)$

Problem

Regression function $\hat{f}(x)$ is discontinuous.

Solution

Weight all points such that their contribution drop off smoothly with distance.

Nadaraya-Watson estimator

$$\hat{f}_{h}(x) = \frac{\sum_{j=1}^{n} K\left(\frac{x-x_{j}}{h}\right) y_{j}}{\sum_{j=1}^{n} K\left(\frac{x-x_{j}}{h}\right)}$$

- *K* is a kernel function where $\int K = 1$, K(a) = K(-a), and $K(0) \ge K(a)$, for all *a*.
- **\hat{f}** is a weighted moving average
- Need to choose *K* and *h*.

Common kernels

Uniform

$$\mathcal{K}(x) = \begin{cases} \frac{1}{2} & -1 < x < 1\\ 0 & \text{otherwise.} \end{cases}$$

Epanechnikov
$$\mathcal{K}(x) = \begin{cases} \frac{3}{4}(1-x^2) & -1 < x < 1\\ 0 & \text{otherwise.} \end{cases}$$

Tri-cube

$$K(x) = \begin{cases} c(1 - |x|^3)^3 & -1 < x < 1 \\ 0 & \text{otherwise.} \end{cases}$$

Gaussian

$$K(x)=\frac{1}{\sqrt{2\pi}}e^{-\frac{1}{2}x^2}$$

Common kernels



KNN vs Smooth Kernel Comparison



Old Faithful

- A smooth kernel is better, but otherwise the choice of kernel makes little difference.
- Optimal kernel (minimizing MSE) is Epanechnikov. It is also fast.
- The choice of *h* is crucial.

Example A

Example B

$$CV(h) = \frac{1}{n} \sum_{j=1}^{n} (y_j - \hat{f}_h^{(-j)}(x_j))^2$$

- (−*j*) indicates *j*th point omitted from estimate
- Pick *h* that minimizes CV.
- Works ok provided there are no duplicate (x, y) pairs. But occasionally odd results.

Many packages available. One of the better ones is KernSmooth:



Problems with the Smooth Weighted Average



Boundary Bias

At some x_0 at a boundary, more of the observations are on one side of the x_0 -The estimated value becomes biased (by those observations).



2 Local polynomials





Kernel smoothing

Equivalent to local constant regression at each point. $\min \sum_{i=1}^{n} w_{i}(x)(y_{i} - a_{0})^{2}$

Local Linear Regression

Fit a line at each point instead.

$$\min \sum_{j=1}^{n} w_j(x)(y_j - a_0 - a_1 x_j)^2$$

Local polynomials



22

$$\min \sum_{j=1}^{n} w_j(x) (y_j - \sum_{k=0}^{p} a_k x_j^p)^2$$

- Local linear and local quadratic are commonly used.
- Robust regression can be used instead
- Less biased at boundaries than kernel smoothing
- Local quadratic less biased at peaks and troughs than local linear or kernel

One useful implementation is

KernSmooth::locpoly(x, y, degree, bandwidth)

- dpill can be used to choose the bandwidth h if degree=1.
- Otherwise, *h* could be selected by cross-validation.
- But most people seem to use trial and error finding the largest h that captures what they think they see by eye.

Best known implementation is loess (locally quadratic)

```
fit <- loess(y ~ x, span=0.75, degree=2,
  family="gaussian", data)</pre>
```

- Uses tri-cube kernel and variable bandwidth.
- span controls bandwidth. Specified in terms of percentage of data covered.
- degree is order of polynomial
- Use family="symmetric" for a robust fit

Old Faithful (Loess, span=0.75)



```
smr <- loess(waiting ~ eruptions, data=faithful)
ggplot(faithful) +
    geom_point(aes(x=eruptions,y=waiting)) +
    ggtitle("Old Faithful (Loess, span=0.75)") +
    geom_line(aes(x=eruptions, y=fitted(smr)),
        col='blue')</pre>
```

Example A (Loess, span=0.75)



Example A (Loess, span=0.22)



Example B (Robust Loess, span=0.75)



30





- Because local polynomials use local linear models, we can easily find standard errors for the fitted values.
- Connected together, these form a pointwise confidence band.
- Automatically produced using geom_smooth

1 Kernel estimators

2 Local polynomials





Splines



34

A spline is a continuous function f(x) interpolating all points (κ_j, y_j) for j = 1, ..., K and consisting of polynomials between each consecutive pair of 'knots' κ_j and κ_{j+1} .

A spline is a continuous function f(x) interpolating all points (κ_j, y_j) for j = 1, ..., K and consisting of polynomials between each consecutive pair of 'knots' κ_j and κ_{j+1} .

- Parameters constrained so that f(x) is continuous.
- Further constraints imposed to give continuous derivatives.
- Cubic splines most common, with f', f'' continuous.

Let $y = f(x) + \varepsilon$ where $\varepsilon \sim IID(0, \sigma^2)$. Then Choose \hat{f} to minimize

$$\frac{1}{n}\sum_{i}(y_i-f(x_i))^2+\lambda\int [f''(x)]^2dx$$

- λ is smoothing parameter to be chosen
- $\int [f''(x)]^2 dx$ is a measure of roughness.
- Solution: f̂ is a cubic spline with knots κ_i = x_i, i = 1, ..., n (ignoring duplicates).
- Other penalties lead to higher order splines
- Cross-validation can be used to select λ .

Smoothing splines



```
smr <- smooth.spline(faithful$eruptions, faithful$waiting,
  cv=TRUE)
smr <- data.frame(x=smr$x,y=smr$y)
ggplot(faithful) +
  geom_point(aes(x=eruptions,y=waiting)) +
  ggtitle("Old Faithful (Smoothing spline, lambda chosen by CV)") +
  geom_line(data=smr, aes(x=x, y=y), col='blue')
```

Smoothing splines



40

Smoothing splines



- Fewer knots than smoothing splines.
- Need to choose the knots rather than a smoothing parameter.
- Can be estimated as a linear model once knots are selected.

- Let $\kappa_1 < \kappa_2 < \cdots < \kappa_K$ be "knots" in interval (a, b).
- Let $x_1 = x$, $x_2 = x^2$, $x_3 = x^3$, $x_j = (x \kappa_{j-3})^3_+$ for $j = 4, \dots, K + 3$.
- Then the regression of *y* on *x*₁,..., *x*_{K+3} is piecewise cubic, but smooth at the knots.
- Choice of knots can be difficult and arbitrary.
- Automatic knot selection algorithms very slow.
- Often use equally spaced knots. Then only need to choose K.

```
fit <- lm(waiting ~ ns(eruptions, df=6), faithful)
ggplot(faithful) +
    geom_point(aes(x=eruptions,y=waiting)) +
    ggtitle("Old Faithful (Natural splines, 6 df)") +
    geom_line(aes(x=eruptions, y=fitted(fit)), col='blue')</pre>
```



Natural splines in R

Example A (Natural splines, 12 df)



Natural splines in R

Example B (Natural splines, 3 df)



Natural splines in R

Example B (Natural splines, 10 df)



47

Splines and geom_smooth()





- Because regression splines use local linear models, we can easily find standard errors for the fitted values.
- Connected together, these form a pointwise confidence band.
- Automatically produced using geom_smooth

1 Kernel estimators

2 Local polynomials





$$y_i = f(\mathbf{x}_i) + \varepsilon_i, \qquad \mathbf{x} \in \mathbb{R}^d$$

Most methods extend naturally to higher dimensions.

- Multivariate kernel methods
- Multivariate local quadratic surfaces
- Thin-plate splines (2-d version of smoothing splines)

$$y_i = f(\mathbf{x}_i) + \varepsilon_i, \qquad \mathbf{x} \in \mathbb{R}^d$$

Most methods extend naturally to higher dimensions.

- Multivariate kernel methods
- Multivariate local quadratic surfaces
- Thin-plate splines (2-d version of smoothing splines)

Problem

The curse of dimensionality!

Curse of dimensionality

Most data lie near the boundary

```
x <- matrix(runif(1e6,-1,1), ncol=100)
boundary <- function(z) { any(abs(z) > 0.95) }
```

```
mean(apply(x[,1,drop=FALSE], 1, boundary))
```

[1] 0.0494

```
mean(apply(x[,1:2], 1, boundary))
```

[1] 0.0971

```
mean(apply(x[,1:5], 1, boundary))
```

[1] 0.2276

Curse of dimensionality

```
Most data lie near the boundary
```

```
x <- matrix(runif(1e6,-1,1), ncol=100)
boundary <- function(z) { any(abs(z) > 0.95) }
```

```
mean(apply(x[,1:10], 1, boundary))
```

```
## [1] 0.4052
```

```
mean(apply(x[,1:50], 1, boundary))
```

[1] 0.9205

```
mean(apply(x[,1:100], 1, boundary))
```

[1] 0.9938

Data are sparse

```
x <- matrix(runif(1e6,-1,1), ncol=100)
nearby <- function(z) { all(abs(z) < 0.5) }
mean(apply(x[,1,drop=FALSE], 1, nearby))</pre>
```

[1] 0.4953

```
mean(apply(x[,1:2], 1, nearby))
```

[1] 0.2512

mean(apply(x[,1:5], 1, nearby))

[1] 0.0291

Data are sparse

```
x <- matrix(runif(1e6,-1,1), ncol=100)
nearby <- function(z) { all(abs(z) < 0.5) }
mean(apply(x[,1:10], 1, nearby))</pre>
```

[1] 5e-04

```
mean(apply(x[,1:50], 1, nearby))
```

[1] 0

```
mean(apply(x[,1:100], 1, nearby))
```

```
## [1] 0
```

lomod <- loess(sr ~ pop15 + ddpi, data=savings)</pre>



library(mgcv)
smod <- gam(sr ~ s(pop15, ddpi), data=savings)</pre>

